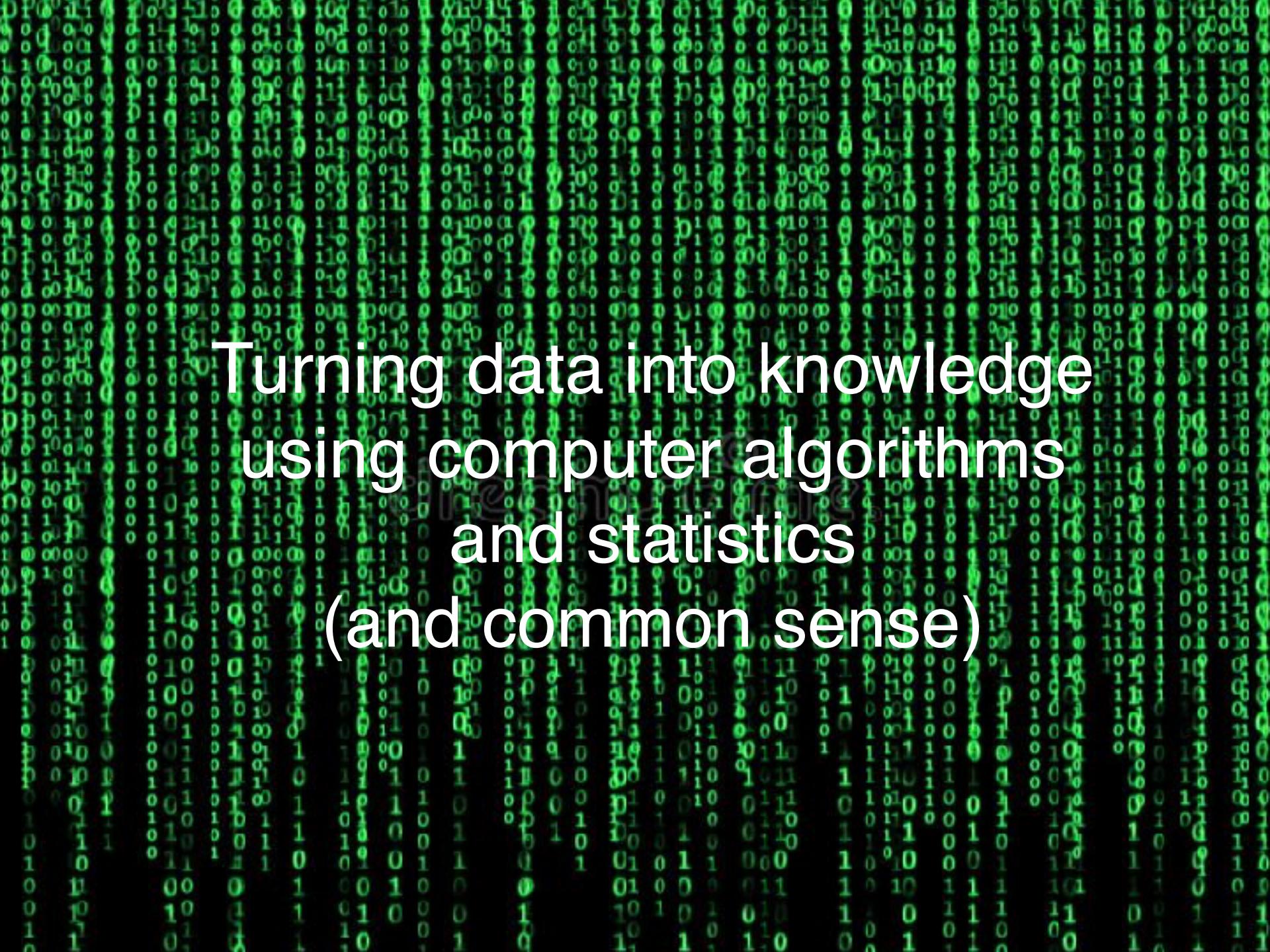


Data Analysis and Machine Learning using Python

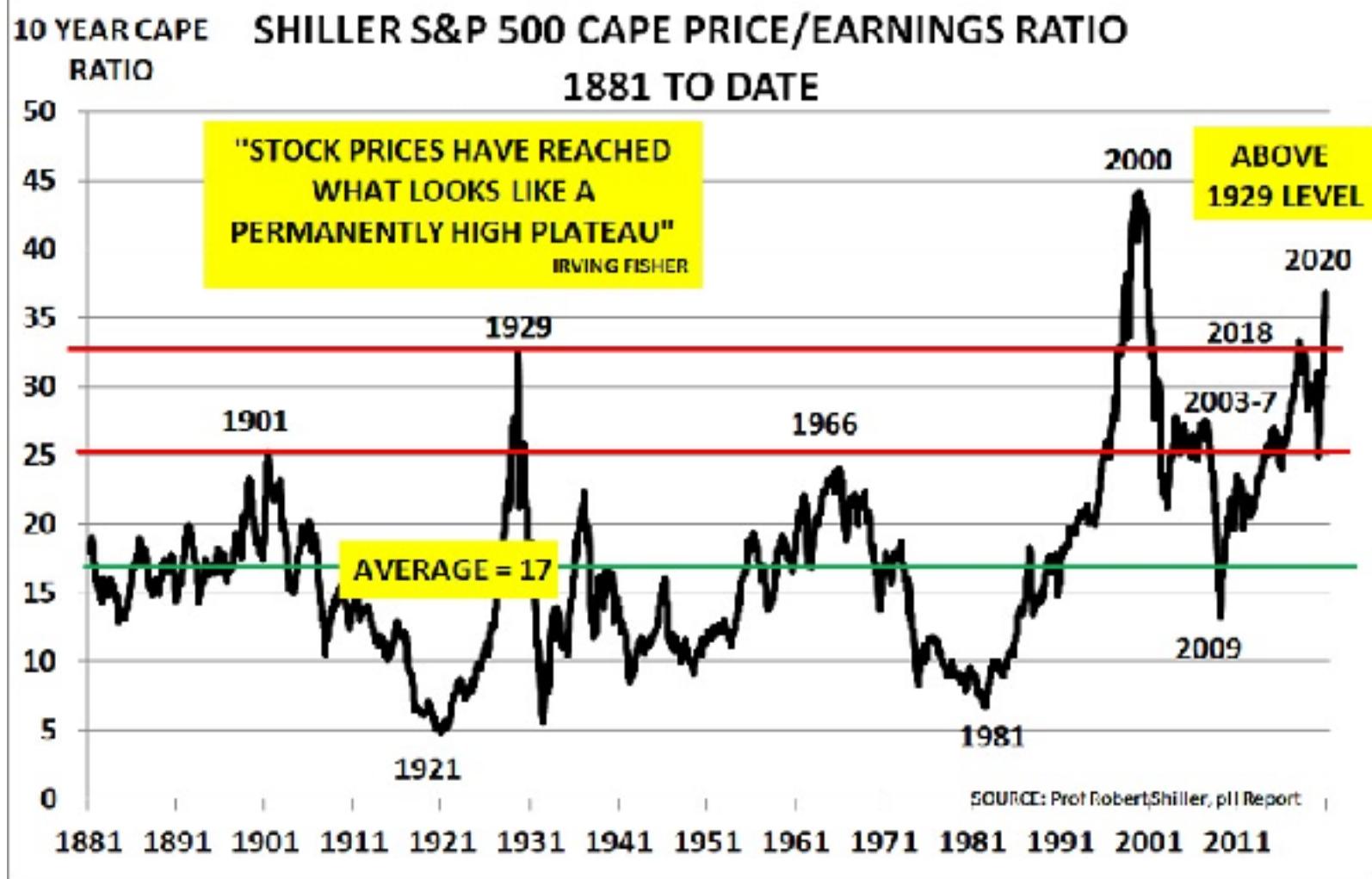
Lecture 1: *Course overview and introduction*
March 9 2024

Gunther Roland
Professor of Physics

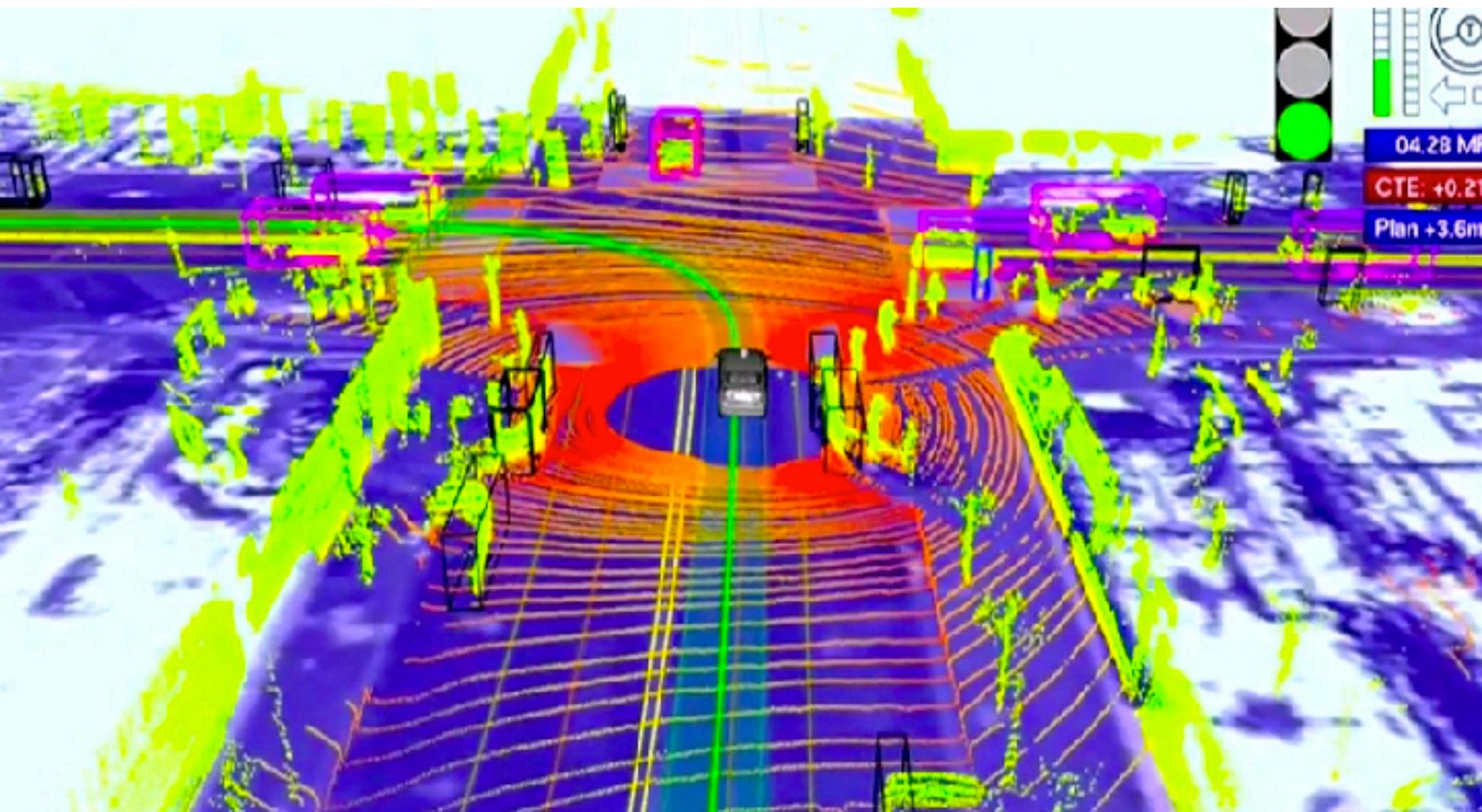




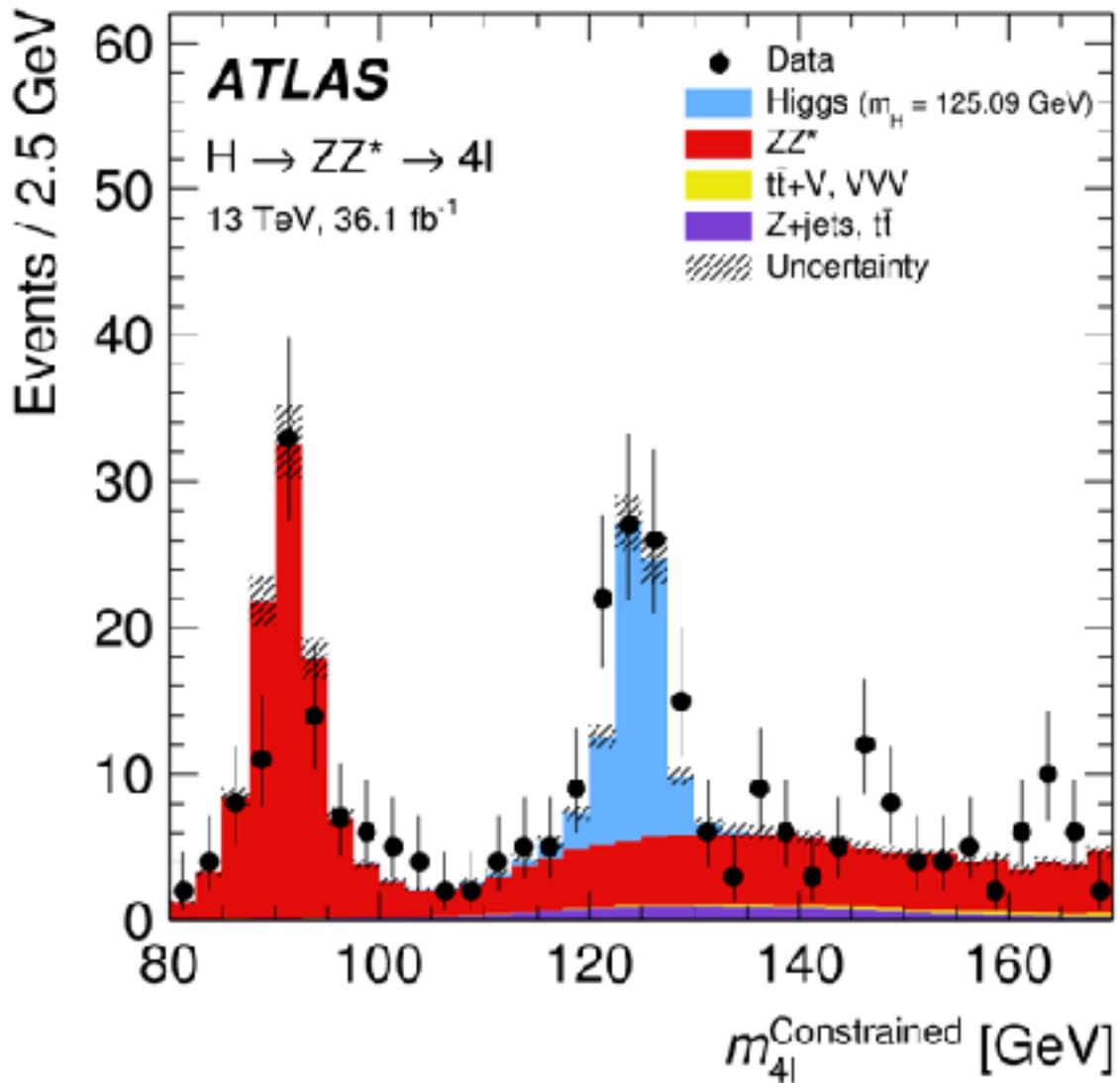
Turning data into knowledge
using computer algorithms
and statistics
(and common sense)



Describe patterns and correlations in financial markets
(and, maybe, predict future)



Make a dynamic model of surroundings
(and, definitely, predict where things will move to)



Discover new laws of nature
(and, possibly, predict other new phenomena)

Me

Prof. Gunther Roland

gunthermroland@gmail.com

Skype: Phobosrolandg

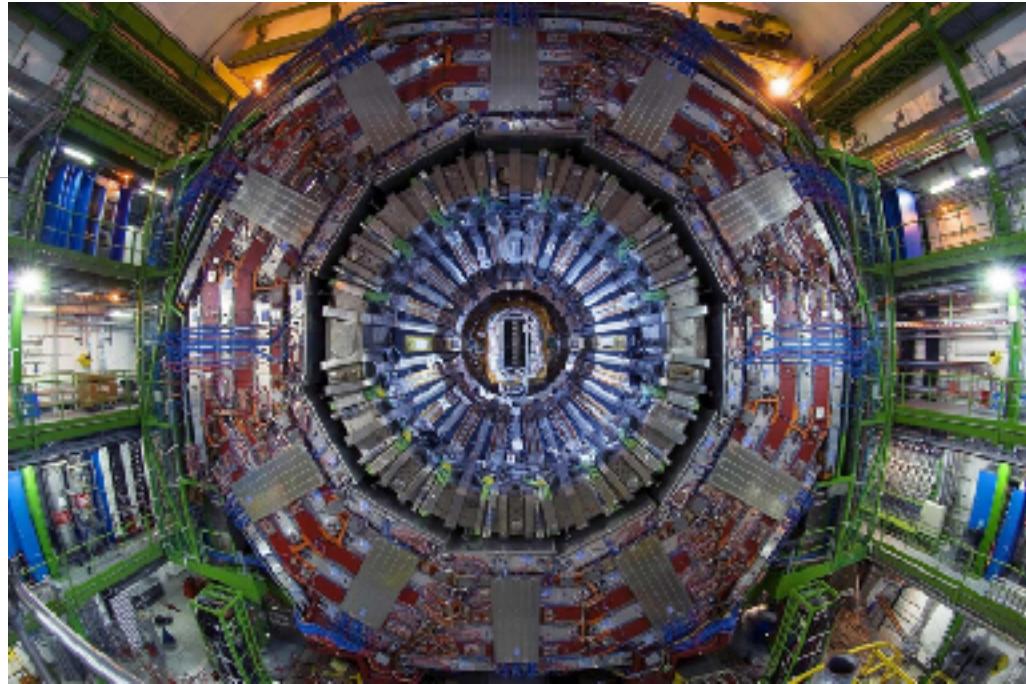
Slack

Experimental physicist

Working at CERN laboratory in Geneva, Switzerland (CMS experiment at LHC accelerator) and BNL laboratory, Long Island, USA (sPHENIX experiment at RHIC accelerator)



Big Data in Physics



- CMS experiment studying subatomic particle collisions
- Studying several billion particle collisions each second
- Produces several Terabyte of data each second → 100 Exabytes/year
 - your laptop probably has ~Terabyte hard drive; this would be 100 million of such hard drives to store
- We are looking for **few** “special” collisions in 100.000.000.000.000.000 to find new laws of nature

What do I want you to learn?

- 10 class sessions
- First 5 sessions: Introduction to Python and basic concepts of data analysis

No.	Topics
Lecture 1	Introduction to Programming in Python <ul style="list-style-type: none">● Variables, Types, Functions● Branching, Iteration, Strings, Lists
Lecture 2	Data Analysis and Basic Statistics <ul style="list-style-type: none">● Scientific Approach to Data Analysis and Basic Statistics● Statistical and Systematic Uncertainties
Lecture 3	Essential Python Libraries for Data analysis <ul style="list-style-type: none">● Using <u>Numpy</u> and <u>Scipy</u>● Plotting using Matplotlib
Lecture 4	Data Visualization and Working with Large Datasets <ul style="list-style-type: none">● Principles of Data Visualization● Best Practices in Working with Large Data Sets● Working with Pandas
Lecture 5	Introduction to Multivariate Analysis <ul style="list-style-type: none">● Introduction to the Concept of Multivariate Analysis● Applications and Significance in Data Analysis

What do I want you to learn?

- Second 5 sessions: Introduction to Machine learning and applications of various ML techniques

Lecture 6	Introduction to Machine Learning <ul style="list-style-type: none">● What is Machine Learning?● Basic Classifications in Machine Learning
Lecture 7	Supervised Learning in Scikit-learn <ul style="list-style-type: none">● Overview and use of Scikit-learn● Decision Tree Algorithm● Multi-layer Perceptrons
Lecture 8	Unsupervised Learning <ul style="list-style-type: none">● K-Means Clustering Algorithm● Principal Component Analysis (PCA)
Lecture 9	Reinforcement Learning <ul style="list-style-type: none">● Concepts and Basic Elements of Reinforcement Learning● Markov Decision Processes
Lecture 10	Deep Learning <ul style="list-style-type: none">● Overview of Deep Learning● Convolutional Neural Networks (CNN)● Recurrent Neural Networks (RNN)● Generative Adversarial Networks (GAN)

Some more thoughts

- Programming background
 - Focus on understanding concepts
 - Focus on essentials, not elegance or **complexity**
 - Not going for the most fancy/short/clever code
- Tools
 - Laptop w/ Python v3.5 or higher
 - Will give examples as Jupyter notebooks

TA sessions

- Please make use of the TA sessions to ask questions and practice
- Also, use the chat feature during the lecture to ask questions
 - the TA will reply right away or let me know so that I can answer questions in the next lecture
- Please let your TA know:
 - What is your programming background?
 - Have you used machine learning techniques before?
 - Do you have Python installed on your laptop?
 - suggest using Jupyter or Spyder
 - Make sure you have Python running for homework

Learning



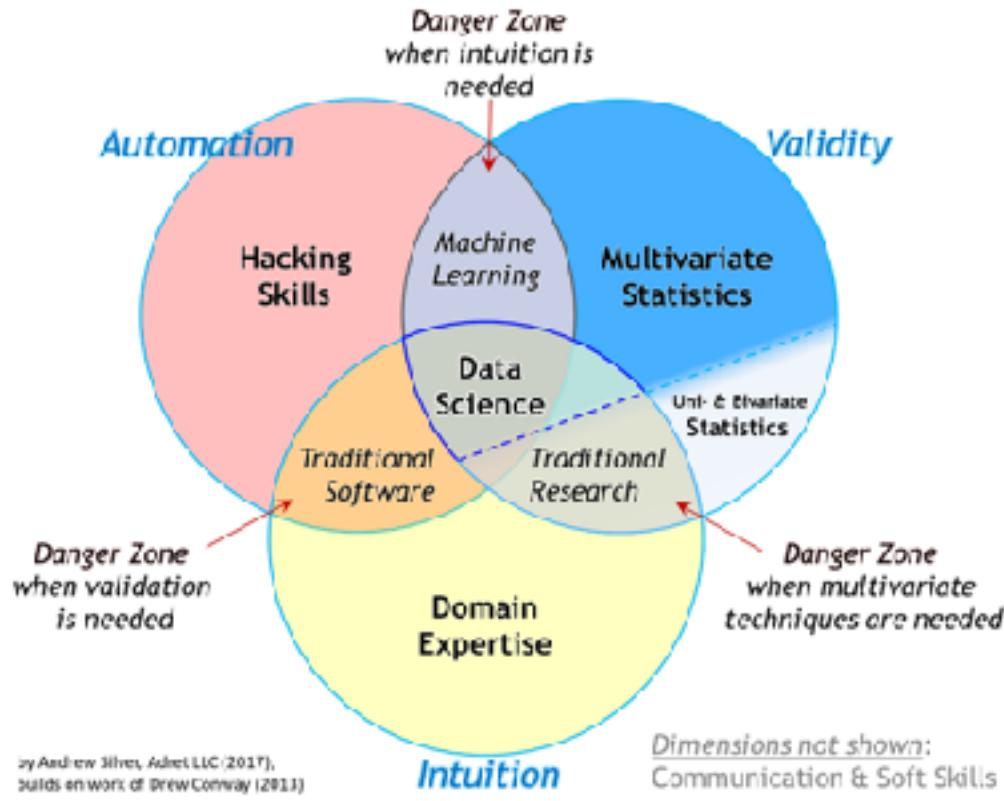
- How many times do you have to watch this clip until you can make the same forehand shot?

Learning by doing



- How many times do you have to watch this clip until you can make the same forehand shot?
- That's not how it works for tennis, and it's the same for any challenging activity - you learn by doing!
- Practice, practice, practice....that's what we are going to do

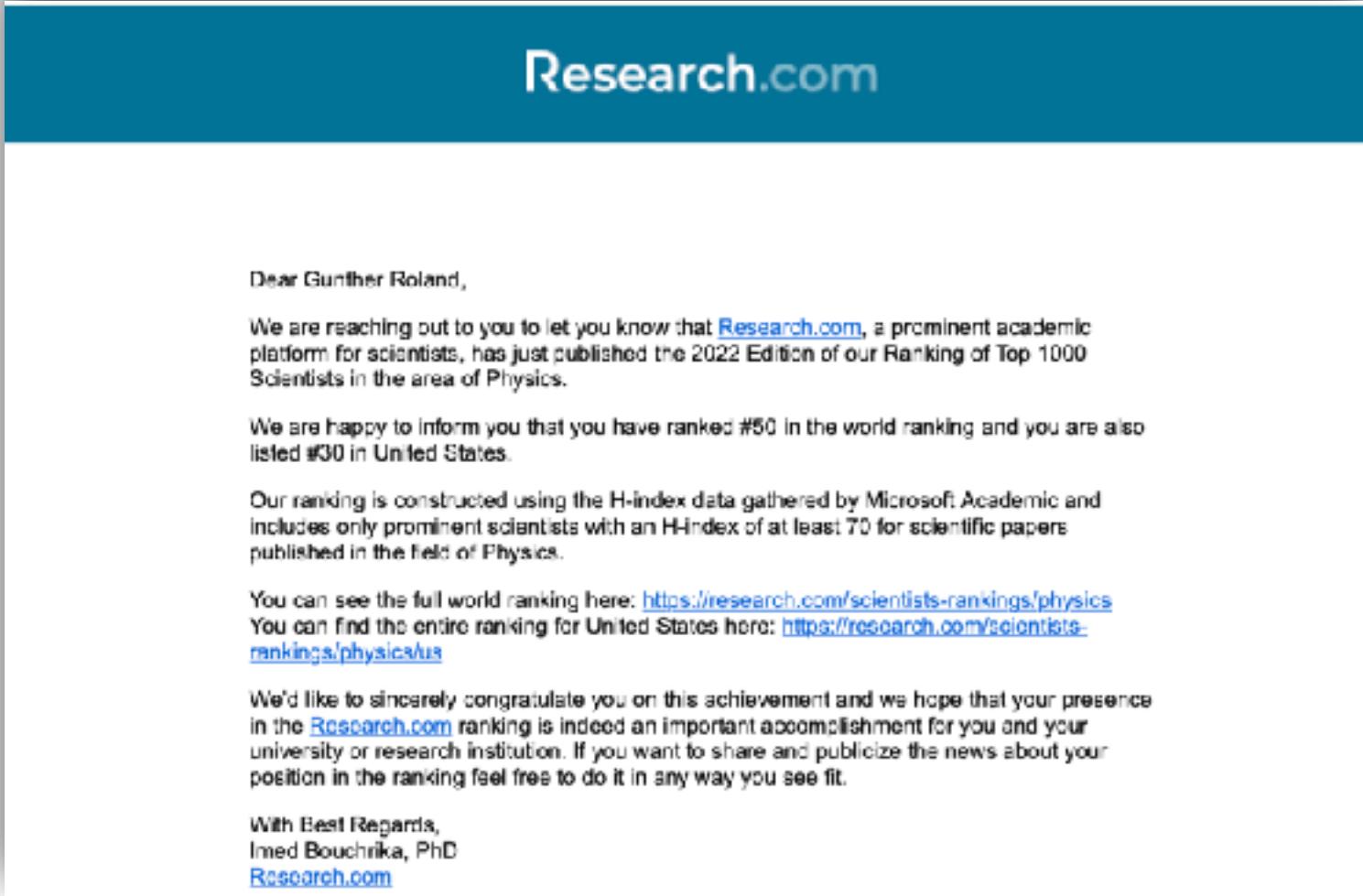
What is Data Science?



wikipedia

Data science is an interdisciplinary field that uses scientific methods, processes, algorithms and systems to extract knowledge and insights from noisy, structured and unstructured data, and apply knowledge from data across a broad range of application domains.

Data science gone wrong



The image shows an email from Research.com. The subject line is "Research.com Ranking of Top 1000 Scientists in Physics". The body of the email starts with a greeting and congratulates the recipient on their ranking. It provides links to the full world ranking and the United States ranking. The email concludes with a closing message and the name of the sender.

Research.com

Dear Gunther Roland,

We are reaching out to you to let you know that [Research.com](#), a prominent academic platform for scientists, has just published the 2022 Edition of our Ranking of Top 1000 Scientists in the area of Physics.

We are happy to inform you that you have ranked #50 in the world ranking and you are also listed #30 in United States.

Our ranking is constructed using the H-index data gathered by Microsoft Academic and includes only prominent scientists with an H-Index of at least 70 for scientific papers published in the field of Physics.

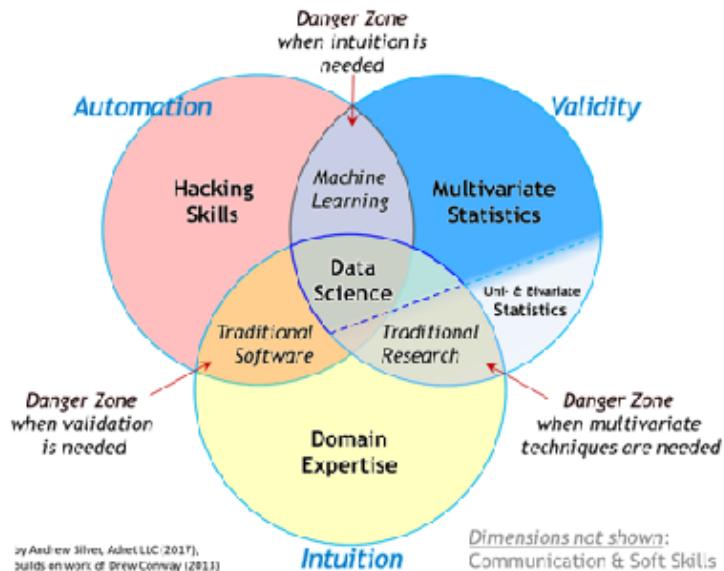
You can see the full world ranking here: <https://research.com/scientists-rankings/physics>
You can find the entire ranking for United States here: <https://research.com/scientists-rankings/physics/us>

We'd like to sincerely congratulate you on this achievement and we hope that your presence in the [Research.com](#) ranking is indeed an important accomplishment for you and your university or research institution. If you want to share and publicize the news about your position in the ranking feel free to do it in any way you see fit.

With Best Regards,
Ined Bouchrika, PhD
[Research.com](#)

email from [research.com](#) earlier this week...

Data science gone wrong



- The “hacking” worked
- The “statistics” are correct
- The conclusions are nonsense
- - did not employ “Domain Expertise”

Top Physics Scientists in United States

This list of top scientists ranking for Physics published on ResearchGate, one of the major websites for scientific research offering cross-models on scientific publications and citations.

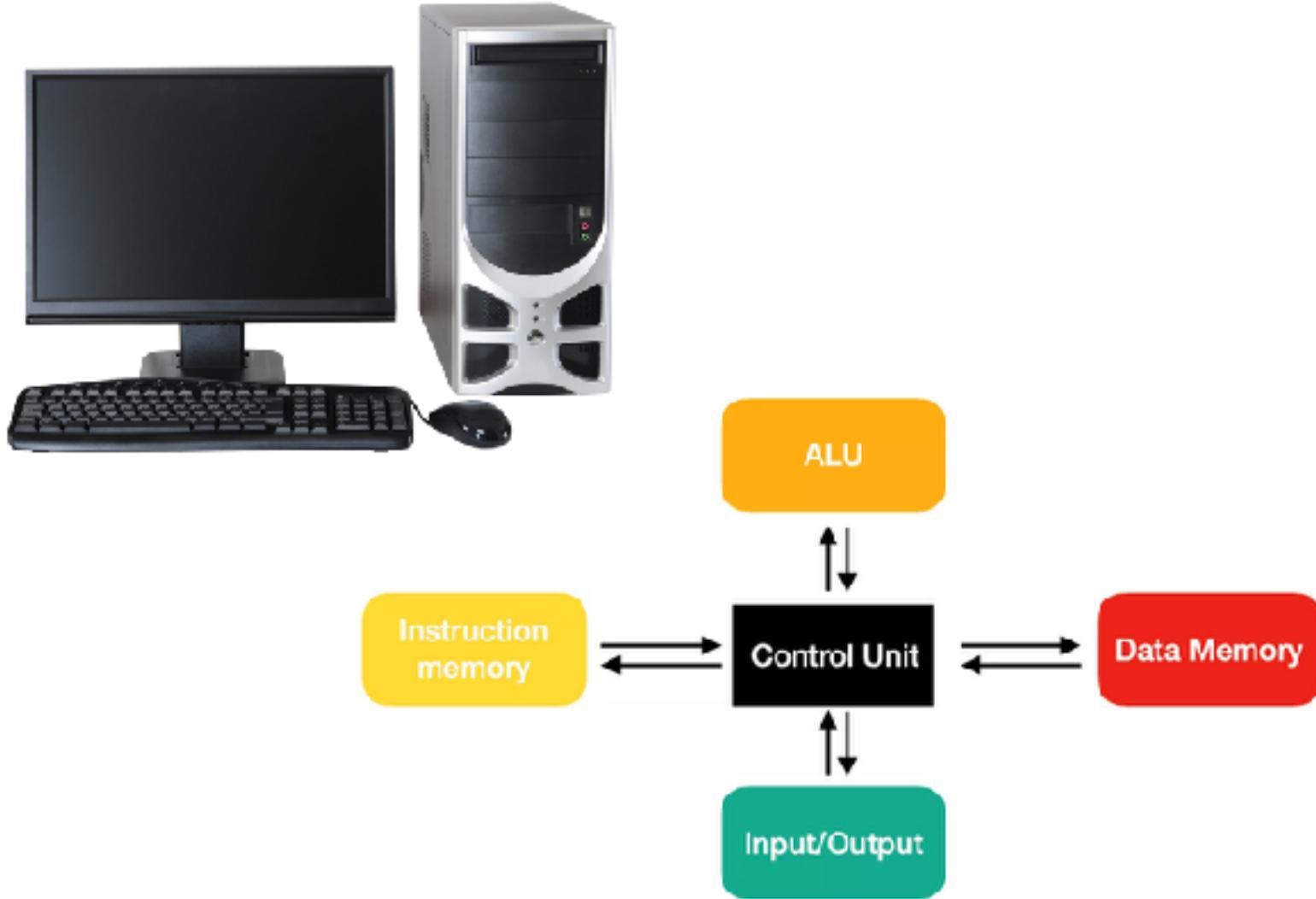
Their h-index, publications and citation counts collected on December 6th 2017. [View here](#)

Scopus indexed		Search	Physics	World	National	Scholar	HIndex	Distance	Publications
1	1		Donald P. Schneider Pennsylvania State University, United States	209	323/73		1,284		
2	2		Michael A. Strauss University of Colorado, United States	262	491/160		1,525		
3	3		Joel Netherland Beamer Permit, United States	251	413/258		2,388		
4	4		J. E. Bjork University of Oregon, United States	223	360/373		1,789		
6	5		Giovanni M. Galeazzi University of Florida, United States	221	361/231		1,817		
11	74		A. V. Krushnai Johns Hopkins University, United States	190	254/276		1,485		
12	27		Julian Borrill Lawrence Berkeley National Laboratory, United States	186	267/347		501		
13	23		Peter Moldrem Johns Hopkins University, United States	186	191/315		1,802		
44	26		T. J. Phelan California Institute of Technology, United States	186	252/260		825		
46	25		Andrew Kingstone University of Bristol, United States	186	182/254		1,617		
47	28		Timothy M. Heckertman Johns Hopkins University, United States	186	107/170		678		
49	29		Bruce A. Bennett Johns Hopkins University, United States	186	129/301		1,256		
50	39		Jennifer Holmes MIT, United States	186	194/288		1,363		

Goals

- Acquire/practice skills for using Python to understand data
- Gain understanding of statistical analysis - what can data tell you, and what is noise?
- Experience with predictive statistics (modeling) and multivariate analysis, including machine learning

What's a computer?



What does it do for us?

- Allows access to data and storage of data
 - 100's of Gigabytes on your laptop
- Performs calculations
 - Billions or trillions of calculations per second
 - Apple M1 CPU allows up to 2.8 TFlops (TFlop = 10^{12} floating point operations per second)
- Built-in calculations (operating system/programming language)
- Complex operations defined by you, the programmer!

Algorithms and programs: think recipe



Moules frites

Prep: 30 mins
Cook: 12 mins
Serves 2



Ingredients

1kg fresh mussels (see tip below)

2 large spring onions

1 large shallot, peeled and halved

1 carrot, peeled and halved lengthways

2 fat garlic cloves, peeled

1 fresh red chilli

1 bunch thyme

Fresh flat-leaf parsley

100ml olive oil

About 150g skinless white fish fillets (Milly Mackerel)

1 tray侯er

2 large dinner forks

For the chips

2 large potatoes, about 300g each, peeled

[preferably Maris Piper, King Edward or Arran]

about 3 drops salt flour

Method

Step 1: Tip the mussels into a large bowl of cold water. Discard any that remain open when tapped; their shells will not shut again. [Note: fresh mussels look black and shiny and should only react periodically to the touch; sea mussels majority should be tightly closed. Avoid any that smell 'fishy', look dry or are mostly open.]

Step 2: Thinly slice the vegetables and garlic. Roll the chilli in your hands to loosen the seeds, then slit it in half and shake out the seeds. Slice the fish into thin stalks, then stalk together and finely chop. Mix over the thyme sprig, discarding the thick stalks, and pick the parsley leaves from the stalks. Prepare the chips (see recipe, below right).

Step 3: Place a large, heavy-based sauté pan (with a lid) on the hob and heat until you can feel a strong heat rising. Throw in the oil, then immediately toss in all of the vegetables, chilli and thyme. The thyme sprig will crackle. If the pan is not enough, cook for about 10 mins, shaking the pan and stirring the vegetables until they start to wilt.

Step 4: With the heat still on high, toss in all the mussels and stir the pan so they form an even layer. Cover with a lid and cook for another 1-2 mins, shaking the pan once or twice.

- Sequence of steps
- Flow control determining order of execution
- A means to stop execution (food is done)
 - Algorithm!
- Programming:
 - translate algorithm into code computer can execute
 - Using specific programming language (e.g., Python)

What is a program?

- Sequence of instructions the computer can execute
- Built from **primitive instructions**, e.g.
 - arithmetic and logic operations
 - simple tests, e.g., for flow control
 - moving data (e.g., input and output, copying data)
- also built from predefined sets of instructions (functions), often provided as **libraries**
- Includes flow control to change sequence depending on input, results and to stop program when done
- Executed (translated into computer's machine code) by special program, i.e, Interpreter or Compiler → **programming language**

Example of an algorithm

- Calculate square root of x , i.e., y such that $y^*y = x$
- Algorithm
 - 1. guess number y
 - 2. if y^*y close enough to x : we're done
 - 3. otherwise average y and x/y to replace y
 - 4. check if y^*y is close enough to x ; otherwise repeat step 3
- Later we'll try to turn this algorithm into a Python program

Our programming language: Python

The screenshot shows the Python.org homepage. At the top, there's a navigation bar with links like "Python", "PSF", "Docs", "PyPI", "Jobs", and "Community". Below the header, the Python logo is prominently displayed next to the word "python". A search bar and a "Donate" button are also visible. The main content area features two code snippets in a terminal window:

```
# Python 3: Simple output (with Unicode)
>>> print("Hello, I'm Python!")
Hello, I'm Python!

# Input, assignment
>>> name = input('What is your name?\n')
>>> print('Hi, ' + name)
What is your name?
Python
Hi, Python.
```

To the right of the code, there's a section titled "Quick & Easy to Learn" with the following text:

Experienced programmers in any other language can pick up Python very quickly, and beginners find the clean syntax and indentation structure easy to learn. [What's your appetite?](#) with our Python 3 overview.

At the bottom, there are five numbered buttons (1, 2, 3, 4, 5) and a summary statement: "Python is a programming language that lets you work quickly and integrate systems more effectively. [» Learn More](#)".

Python is a programming language that lets you work quickly
and integrate systems more effectively. [» Learn More](#)

Python

Python is a [high-level](#), [interpreted](#), [general-purpose programming language](#). Its design philosophy emphasizes [code readability](#) with the use of [significant indentation](#).^[31]

Python is [dynamically-typed](#) and [garbage-collected](#). It supports multiple [programming paradigms](#), including [structured](#) (particularly [procedural](#)), [object-oriented](#) and [functional programming](#). It is often described as a "batteries included" language due to its comprehensive standard library.^{[32][33]}

We'll discuss what these words mean

Wikipedia

- Invented in the Netherlands by Guido van Rossum (early 90's)
- Named after Monty **Python**
- Initially often thought of as *scripting language*
 - now used for much more
- Open source
- Extensive selection and use of open-source *libraries*
- One of the most popular programming languages

Visualizing data with plots is essential: Matplotlib

The screenshot shows the homepage of matplotlib.org. At the top, there's a navigation bar with links to various organizations and a search bar. Below the header, the **matplotlib** logo is on the left, and a menu bar with links to **Plot types**, **Examples**, **Tutorials**, **Reference**, **Usage guide**, **Develop**, and **Release notes**. To the right of the menu are social media icons for LinkedIn, YouTube, GitHub, and Twitter.

The main content area features a title **Matplotlib: Visualization with Python**. Below it is a brief description: **Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python. Matplotlib makes easy things easy and hard things possible.** To the left of the text is a small plot showing several blue barbs (arrows) representing vector fields, with the caption **barbs(X, Y, U, V)**.

Below the description is a bulleted list of features:

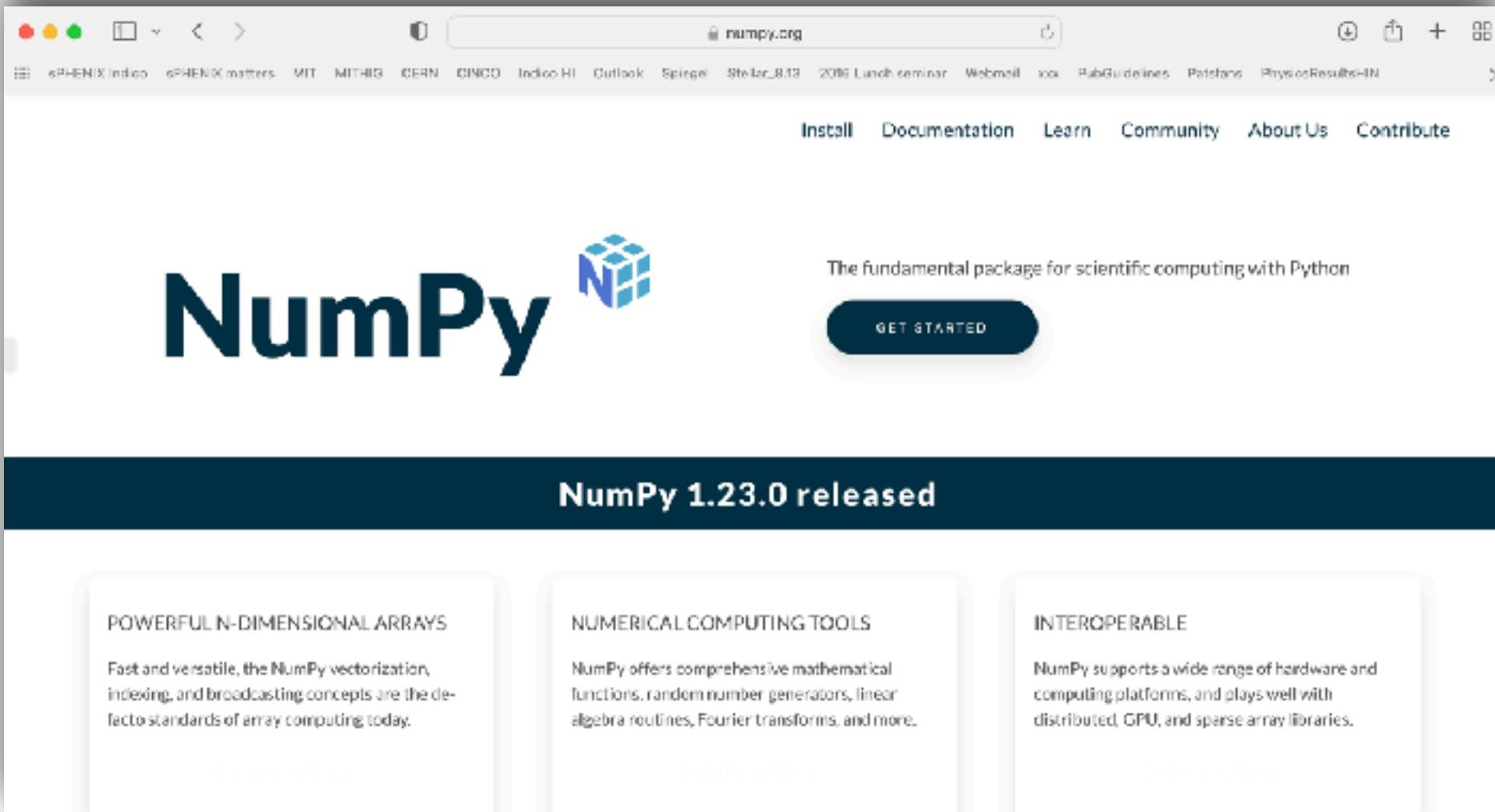
- Create publication quality plots.
- Make [interactive figures](#) that can zoom, pan, update.
- Customize visual style and layout.
- Export to many file formats .
- Embed in [JupyterLab](#) and [Graphical User Interfaces](#).
- Use a rich array of [third-party packages](#) built on Matplotlib.

At the bottom of the main content area is a purple button with white text: **Try Matplotlib (on Binder)** followed by a right-pointing arrow.

Below the main content are five navigation links with icons:

- Getting Started** (monitor icon)
- Examples** (file folder icon)
- Reference** (book icon)
- Cheat Sheets** (document icon)
- Documentation** (colorful circular icon)

Need a toolbox of numerical algorithms: NumPy



The screenshot shows the NumPy website at numpy.org. The page features a large "NumPy" logo with a blue cube icon. Below the logo, a banner announces "NumPy 1.23.0 released". Three main features are highlighted in boxes: "POWERFUL N-DIMENSIONAL ARRAYS", "NUMERICAL COMPUTING TOOLS", and "INTEROPERABLE". A "GET STARTED" button is visible above the release banner.

The NumPy logo consists of the word "NumPy" in a bold, dark blue sans-serif font, followed by a small blue 3D cube icon where the letter "N" would be.

The banner text "NumPy 1.23.0 released" is centered in a white box with a dark blue background.

The three main features are:

- POWERFUL N-DIMENSIONAL ARRAYS**: Described as "Fast and versatile, the NumPy vectorization, indexing, and broadcasting concepts are the de-facto standards of array computing today."
- NUMERICAL COMPUTING TOOLS**: Described as "NumPy offers comprehensive mathematical functions, random number generators, linear algebra routines, Fourier transforms, and more."
- INTEROPERABLE**: Described as "NumPy supports a wide range of hardware and computing platforms, and plays well with distributed, GPU, and sparse array libraries."

A "GET STARTED" button is located below the banner.

and a math and science toolkit: SciPy

The screenshot shows the SciPy.org homepage in a web browser. The address bar displays "scipy.org". The top navigation bar includes links for "Install", "Documentation", "Download", "Community", "About Us", and "Contribute". Below the header, the SciPy logo is prominently displayed. To the right of the logo is the text "Fundamental algorithms for scientific computing in Python" and a "GET STARTED" button. A dark banner at the bottom of the page announces "Join us at SciPy 2022 July 11th - 17th!" followed by the date "2022-06-27". The page is divided into three main sections: "FUNDAMENTAL ALGORITHMS", "BROADLY APPLICABLE", and "FOUNDATIONAL", each with descriptive text.

FUNDAMENTAL ALGORITHMS
SciPy provides algorithms for optimization, integration, interpolation, eigenvalue problems, algebraic equations, differential equations, statistics and many other classes of problems.

BROADLY APPLICABLE
The algorithms and data structures provided by SciPy are broadly applicable across domains.

FOUNDATIONAL
Extends NumPy providing additional tools for array computing and provides specialized data structures, such as sparse matrices and k-dimensional trees.

and for machine learning: scikit-learn

The screenshot shows the official scikit-learn website at scikit-learn.org/stable/. The page features a large header with the scikit-learn logo and navigation links for Install, User Guide, API, Examples, Community, and More. Below the header, there's a main title "scikit-learn" and subtitle "Machine Learning in Python". A prominent orange sidebar on the right lists key features: Simple and efficient tools for predictive data analysis, Accessible to everybody, and reusable in various contexts, Built on NumPy, SciPy, and matplotlib, and Open source, commercially usable - BSD license. The main content area is divided into three sections: Classification, Regression, and Clustering, each with a brief description, applications, algorithms, and a corresponding visualization. The Classification section shows a 3x5 grid of small plots illustrating various classification models. The Regression section shows a line plot titled "Boosted Decision Tree Regression" comparing training samples and estimators. The Clustering section shows a scatter plot of digits data points colored by cluster assignment.

scikit-learn
Machine Learning in Python

Getting Started Release Highlights for 1.1 GitHub

Classification

Identifying which category an object belongs to.

Applications: Spam detection, image recognition.

Algorithms: SVM, nearest neighbors, random forest, and more...

Regression

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, Stock prices.

Algorithms: CVR, nearest neighbors, random forest, and more...

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes.

Algorithms: k-Means, spectral clustering, mean-shift, and more...

Examples

Examples

Examples

Running Python

- Which computer/operating system are you using (Windows, Mac OSX, Linux?)
- Do you have Python installed?
- Recommend use of an IDE (e.g., spyder, Jupyter, etc)

Proposed programming environment (IDE): jupyter

jupyter triangle (autosaved)  Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel) O

In [3]:

```
from scipy.optimize import curve_fit
import matplotlib.pyplot as plt
import skimage
import numpy as np

xmin = -1
xmax = 1

# Generate Gaussian data points
eta1 = np.random.uniform(xmin,xmax,1000)
eta2 = np.random.uniform(xmin,xmax,1000)

for i in range(10):
    print(eta1[i], eta2[i])

data = []
for i in range(len(eta1)):
    for i in range(len(eta2)):
```

For installing jupyter, I use ANACONDA

The screenshot shows a web browser window displaying the Anaconda website at www.anaconda.com. The page features a large green header with the word "ANACONDA". Below the header, there's a navigation bar with links for "Products", "Pricing", "Solutions", "Resources", "Partners", "Blog", and "Company". A "Contact Sales" button is also visible. The main content area has a white background with a large green title: "Data science technology for a better world." Below the title, a green paragraph explains that Anaconda offers the easiest way to perform Python/R data science and machine learning on a single machine, mentioning thousands of open-source packages and libraries. A prominent green "Download" button with a white arrow icon is centered below this text. Below the button, smaller text indicates the download is for Mac OS, is a 64-bit graphical installer, and is 591 MB in size. There's also a link to "Get Additional Installers". At the bottom of the page, a dark footer bar contains the text "This website uses cookies to ensure you get the best experience on our website. [Privacy Policy](#)" and a "Accept" button.

and we'll communicate over Slack

The screenshot shows the Slack website homepage with a dark purple background. At the top, there's a navigation bar with links for Product, Solutions, Enterprise, Resources, Pricing, Sign In, TALK TO SALES, and TRY FOR FREE. Below the navigation, a large white text area contains the headline "Great teamwork starts with a digital HQ". Underneath the headline is a subtext: "With all your people, tools and communication in one place, you can work faster and more flexibly than ever before." There are two sign-up buttons: "SIGN UP WITH EMAIL" and "SIGN UP WITH GOOGLE". To the right of the text area, there's a graphic showing a smartphone displaying a Slack channel named "#projectunicorn" and a laptop screen showing a similar interface. At the bottom, a section titled "TRUSTED BY COMPANIES ALL OVER THE WORLD" lists logos for Airbnb, Uber, Spotify, Netflix, One Medical, TD Ameritrade, Intuit, NASA, and TIME.

Slack

- How many of you have an account on Slack already?
- Let's take a 5min break and see if you can create your Slack account, start downloading Anaconda, if using, (and get some refreshments, take a bio-break etc)

Terminology

- Programming language provides set of primitive operations (primitives), e.g., “=”, “+”
- Primitives can be combined (with certain rules) into expressions, e.g. “ $a = b+x$ ”
- Expressions manipulate data and can themselves have a value, e.g., “ $1+1 == 2$ ” has Boolean value “True”
- Expressions can be typed in Python shell directly one-by-one or stored in a file and executed in sequence → program (“Python script”)

A Python script

```
1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Created on Sun Jul 10 15:47:37 2022
5
6  @author: rolandg
7  """
8
9  # say hello (this is a comment)
10 print("Hello world!")
11
12 # that's the end of my script
13 #
```

- this is the directory where my Python scripts are stored
- when running the script (e.g. in spyder):

```
In [183]: runfile('/Users/rolandg/Python/Hello.py', wdir='/Users/rolandg/Python')
Hello world!
```

```
In [184]: |
```

this is the directory where my Python scripts are stored



Data

- Expressions can be used to manipulate data
- Data in Python are **data objects**
- **Data objects** have different **types**
- **Type** describes nature of data (e.g. string of characters or floating point number or complex data structure)
- **Type** determines the kind of operations possible on data structure
- Data objects can be scalar (no internal structure, e.g., an integer number) or non-scalar (internal structure, e.g., a list of integer numbers)

Simple data types

```
In [140]: i = 5
```

```
In [141]: type(i)
```

```
Out[141]: int
```

```
In [142]: x = 25.5
```

```
In [143]: type(x)
```

```
Out[143]: float
```

```
In [144]: a = "Gunther"
```

```
In [145]: type(a)
```

```
Out[145]: str
```

```
In [146]: f = True
```

```
In [147]: type(f)
```

```
Out[147]: bool
```

- Unlike e.g., C++ (and most other compiled languages) types in Python are usually not assigned explicitly, but determined by Python from context/data/assignment
- usually convenient, but can lead to surprises

Types can be converted

```
In [148]: i = 100
```

```
In [149]: x = float(i)
```

```
In [150]: print(x)  
100.0
```

```
In [151]: type(x)  
Out[151]: float
```

Combine objects and operators to form expressions

```
In [152]: x = 100.0  
In [153]: y = 10.0  
In [154]: z = x/y  
In [155]: print(z)  
10.0
```

Algebraic, e.g., +,-,*,/

```
In [159]: y == z  
Out[159]: True
```

```
In [160]: y = x  
Out[160]: False
```

```
In [161]: y < x  
Out[161]: True
```

Logic, e.g., ==, <, >

note difference
between “=” and “==”

Expressions (“x/z” or “y == x”) have value and type

Expressions

- Type of expression is usually determined from input object type

```
In [162]: i = 10 + 10
```

```
In [163]: type(i)  
Out[163]: int
```

```
In [164]: x = 10.0 + 10.0
```

```
In [165]: type(x)  
Out[165]: float
```

- Parenthesis can be used to fix order of computation

```
In [170]: x = 10.0 + 20.0/2.0
```

```
In [171]: print(x)  
20.0
```

```
In [172]: x = (10.0 + 20.0)/2.0
```

```
In [173]: print(x)  
15.0
```

Equal sign “=”

- Equal sign assigns value to a variable (object)

```
In [180]: e_approx = 2.718
```

```
In [181]: print(e_approx)  
2.718
```

```
In [182]: log(e_approx)  
Out[182]: 0.999896315728952
```

- assignment binds name (e_approx) to value (2.718) stored in computer memory
- can bind name to another value with a new assignment
 - at that point, old value can't be accessed any more

Some basics for understanding Python code

Whitespace is meaningful in Python: especially **indentation** and placement of newlines

- Use a newline to end a line of code
 - Use \ when must go to next line prematurely
- No braces {} (as in C++) to mark blocks of code, use *consistent* indentation instead
 - First line with *less* indentation is outside of the block
 - First line with *more* indentation starts a nested block
- Colons start of a new block in many constructs, e.g. function definitions

Some basics for understanding Python code

- **Indentation matters (unlike e.g., C++)**
 - Block structure indicated by indentation

```
x = 9.0
print("With indentation")
if(x > 10.0):
    y = 10
    print("Yeah")

print("Without indentation")
if(x > 10.0):
    y = 10
print("Yeah")
```

```
In [185]: %runfile('/Users/rolandg/Python/Hello.py', wdir='/Users/rolandg/Python')
Hello world!
With indentation
Without indentation
Yeah
```

Some basics for understanding Python code

- First assignment to a variable (name) creates it
 - Variable types don't need to be declared
 - Python figures out the variable types automatically
- Assignment is `=` and comparison is `==`
- For numbers `+ - * / %` work as expected
 - For strings, `+` provides concatenation
- Logical operators are words: “`and`”, “`or`”, “`not`”
- The basic printing command is `print`

Operators do different things for different types

```
In [186]: i = 10
In [187]: j = 20
In [188]: print(i+j)
30

In [189]: a = "abc"
In [190]: b = "def"
In [191]: print(a+b)
abcdef
```

Variable assignment

- *Binding a variable* means setting a *name* to hold a *reference* to some *object or value* in memory
 - Assignment creates references, not copies
- Names in Python do not have an intrinsic type, objects have types
 - Python determines the type associated with named object automatically, based on what data is assigned to it
- You create a name the first time it appears on the left side of an assignment expression:
`x = 3`
- A reference is deleted via garbage collection after any names bound to it have passed out of scope

Naming conventions

- Names are **case sensitive** and cannot start with a number. They can contain letters, numbers, and underscores, e.g.,
bob Bob _bob _2_bob_ bob_2 BoB
 - I do NOT recommend using upper/lower case to distinguish variables (e.g., bob vs Bob) - just asking for confusion
- There is a number of reserved words that you can't re-use or redefine as they are intrinsic Python commands:
`and, assert, break, class, continue,
def, del, elif, else, except, exec,
finally, for, from, global, if, import,
in, is, lambda, not, or, pass, print,
raise, return, try, while`

Naming conventions

- Recommended naming conventions:
- **joined_lower** for functions, methods and attributes
- **ALL_CAPS** for constants
- **StudlyCaps** for classes
- **camelCase** if needed to conform to pre-existing conventions
- Sticking to these conventions will make your code easier to read for others (and for yourself)
- For temporary variables I usually use i,j,k... for integers, x,y,z,.. for floats and a,b,c... for strings

Some options for assignments

- Make multiple assignments at the same time

```
In [18]: x,y = 4,8
```

```
In [19]: print(x,y)  
4 8
```

Can swap variables very easily

```
In [20]: x,y = y,x
```

```
In [21]: print(x,y)  
8 4
```

- Assignments can be chained (but not sure if I recommend overusing this)

```
In [22]: x = y = z = 123
```

```
In [23]: print(x,y,z)  
123 123 123
```

However, I suggest that for the beginning we write things as explicitly as we can - make it easy to read for non-Python experts

Control flow and branching

- For almost any program to be useful we need to control the flow of commands based on data and conditions we encounter → control flow, branching and iteration

Branching: “if”, “else”, “elif”

- Decide whether block of code should be executed based on a Boolean condition (“True” or “False”)

```
if <condition>:  
    <expression>  
    <expression>  
    ...
```

```
if <condition>:  
    <expression>  
    <expression>  
    ...  
else:  
    <expression>  
    <expression>  
    ...
```

```
if <condition>:  
    <expression>  
    <expression>  
    ...  
elif <condition>:  
    <expression>  
    <expression>  
    ...  
else:  
    <expression>  
    <expression>  
    ...
```

- Recall importance of indentation; don’t forget “.”

Iteration: “while” loops

```
while <condition>:  
    <expression>  
    <expression>  
    ...
```

- Enter loop if <condition> is “True”
- Continue executing expressions in loop repeatedly until condition becomes “False”
- If condition remains “True” will loop forever...

Iteration: “for” loops

- Iterate through a range of numbers

```
for <variable> in range(<some_num>):  
    <expression>  
    <expression>  
    ...
```

- Start with variable in smallest value of range
- Execute block of code
- Increase variable by 1 and check if still in range
 - if True, execute block of code
 - if False, exit for loop
- Can do the same with a while loop, but needs more lines of code and easier to make mistakes

The “range” function in Python

- *range(start, stop, step)*
- Default values are *start = 0* and *step = 1*
- will count from *0* to *stop-1* in steps of *1*
- can change start and stop values, as well as step size

“Sequence types”

- What if you have many numbers, letters etc that you want to access?
- → Sequence types: Tuples, strings and lists
- Can be “mutable” and “immutable”
- Lists and tuples can contain objects of different types
 - although in most cases we will use lists containing objects of uniform type

Examples

- **Tuple:** ('Gunther', 'abc', 123, [SomeOtherSequence])
 - An *immutable* sequence of items
 - Items can be of mixed types,
 - can include collection types
 - note the parenthesis
- **Strings:** "Gunther Roland"
 - *Immutable* sequence of characters
- **List:** [1, 2, 3, 'Gunther', ('left', 'right')]
 - *Mutable* sequence of items of mixed types
 - note the square brackets

Sequence type syntax

- Tuples, strings and list share a lot of functionality and syntax (i.e., the way to access elements)
- Key difference:
 - lists are *mutable* (items in the list can be changed),
 - tuples and strings are immutable (e.g., cannot change the “o” in “Roland” to a different letter) - *immutable*
-

Sequence type creation - syntax

- Create tuples using parentheses and commas

```
In [24]: tu = (23,45,"abc","def",54)  
In [25]: print(tu)  
(23, 45, 'abc', 'def', 54)
```

- In [26]: type(tu)
Out[26]: tuple

- Create lists are using square brackets and commas

```
In [27]: li = [23,45,"abc","def",54]  
In [28]: print(li)  
[23, 45, 'abc', 'def', 54]  
In [29]: type(li)  
Out[29]: list
```

- Create strings using quotes (" , ' , or """")

```
In [31]: st = "23,45,abc,def,54"  
In [32]: print(st)  
23,45,abc,def,54  
In [33]: type(st)  
Out[33]: str
```

Accessing member elements of sequences

- Use square bracket “array” notation

```
In [34]: print(tu[1])
```

```
45
```

```
In [35]: print(li[1])
```

```
45
```

```
In [36]: print(st[1])
```

```
3
```

Note that we are counting **from 0** - first element is accessed by li[0]

Interesting feature in Python (for a C++ programmer)

```
In [37]: print(li[0])
```

```
23
```

```
In [38]: print(li[-1])
```

```
54
```

```
In [39]: print(li[-5])
```

```
23
```

Positive index: count from “left”, starting with 0

Negative index: count from “right”, starting with -1

Important operation: “slicing”

- Return a sequence of same type, containing a copy of a subset of the original sequence

```
In [40]: li = [23, 45, "abc", "def", 54]
```

```
In [41]: print(li[1:3])
[45, 'abc']
```

```
In [42]: type(li[1:3])
Out[42]: list
```

- Start copying at the first index
- Stop copying **before** second index
- Negative indices count from “right”

Copying the whole sequence

- [:] makes a *copy* of the entire sequence
- IMPORTANT: Note the difference between these two operations for lists:

```
In [43]: li = [23,45, "abc", "def", 54]
```

```
In [44]: li2 = li[:]
```

```
In [45]: li3 = li
```

```
In [46]: li[0] = 11
```

```
In [47]: print(li)
[11, 45, 'abc', 'def', 54]
```

```
In [48]: print(li2)
[23, 45, 'abc', 'def', 54]
```

```
In [49]: print(li3)
[11, 45, 'abc', 'def', 54]
```

Makes a copy

Points to same list

Useful operator: “in”

- Test whether a value is inside sequence (True or False):

```
In [50]: tu = (1,2,3,5,6)
```

```
In [51]: 1 in tu  
Out[51]: True
```

```
In [52]: 2 in tu  
Out[52]: True
```

- In [53]: 4 in tu
Out[53]: False

- Tests whether substrings is inside string:

```
In [54]: st = "Gunther"
```

```
In [55]: "un" in st  
Out[55]: True
```

```
In [56]: "e" in st  
Out[56]: True
```

- In [57]: "a" in st
Out[57]: False

Useful operator: “+”

The “+” operator produces a *new* sequence by concatenation of arguments:

```
In [58]: (2,4,6)+(1,2,3)
Out[58]: (2, 4, 6, 1, 2, 3)
```

```
In [59]: "Gunther" + "Roland"
Out[59]: 'GuntherRoland'
```

Somewhat useful operator: “*”

The “*” operator produces a *new* sequence by concatenating arguments n times:

```
In [60]: "Gunther" * 4
Out[60]: 'GuntherGuntherGuntherGunther'
```

```
In [61]: (1,2,3) * 4
Out[61]: (1, 2, 3, 1, 2, 3, 1, 2, 3)
```

“Mutability” - lists are mutable

```
In [64]: li = [1,2,3,4,5]
In [65]: li[4] = li[0]
In [66]: print(li)
[1, 2, 3, 4, 1]
```

- Changes list *in place* - change content of memory list is pointing to

“Mutability” - tuples are NOT mutable

```
In [67]: tu = (1,2,3,4,5)
In [68]: tu[4] = tu[0]
Traceback (most recent call last):
  File "<ipython-input-68-f2c37bc18289>", line 1, in <module>
    tu[4] = tu[0]
TypeError: 'tuple' object does not support item assignment
```

- Cannot change a tuple in place
- If you want a different value in tuple, need to make new one
- Why tuples? Immutability means operations are faster on tuples than on lists

Important operators for lists: “append” and “insert”

Lists are mutable - can be extended:

```
>>> li = [1, 11, 3, 4, 5]
```

```
>>> li.append('a')    # Note the method  
syntax
```

```
>>> li  
[1, 11, 3, 4, 5, 'a']
```

```
>>> li.insert(2, 'i')
```

```
>>> li  
[1, 11, 'i', 3, 4, 5, 'a']
```

Other powerful built-in functions for lists

Lists have many methods, including `index`, `count`, `remove`, `reverse`, `sort`

```
>>> li = ['a', 'b', 'c', 'b']
>>> li.index('b')      # index of 1st occurrence
1
>>> li.count('b')     # number of occurrences
2
>>> li.remove('b')    # remove 1st occurrence
>>> li
['a', 'c', 'b']
```

Other powerful built-in functions for lists

```
>>> li = [5, 2, 6, 8]

>>> li.reverse()      # reverse the list *in place*
>>> li
[8, 6, 2, 5]

>>> li.sort()         # sort the list *in place*
>>> li
[2, 5, 6, 8]

>>> li.sort(some_function)
# sort in place using user-defined comparison
```

More on tuples

- The **comma** is the tuple creation operator, not parenthesis

```
>>> 1,  
(1,)
```

- Parenthesis is just shown for clarity

```
>>> (1,)  
(1,)
```

- No comma, no tuple

```
>>> (1)  
1
```

- Empty tuples have a special syntactic form

```
>>> ()  
()  
>>> tuple()  
()
```

tuples vs lists

- Lists slower but more powerful than tuples
 - Lists can be modified, and they have lots of handy operations and methods
 - Tuples are immutable and have fewer features
- Can convert between tuples and lists using the list() and tuple() functions:

```
li = list(tu)
tu = tuple(li)
```
- I find myself using lists more often - can do anything tuples can

Input

- Often (always?) need to get input data from user or from file on disk
- Asking for user input: “`input()`” (what else?)

```
In [7]: a = input("enter a name:")
```

```
enter a name:Gunther
```

```
In [8]: print(a)
```

```
Gunther
```

```
In [9]: |
```

Input from text file

- We will mostly use text files as data source
- Numbers or text entries, separated by “,” (“CSV”), whitespace “ “ or some other delineator
- e.g.:

```
14.Sep.B01.01,1,9,0.2,10.4,0.2,36.8,0.3,3400,200
14.Sep.B01.02,2,9.1,0.6,10.5,1,36.8,0.8,3500,400
14.Sep.B01.03,3,8.9,0.1,10.1,0.3,36.5,0.3,3281,0.44
14.Sep.B01.04,4,9,0.5,10.4,0.5,36.9,0.5,3400,0.9
14.Sep.B01.05,5,8.9,0.1,10.4,0.1,37,0.1,3400,51
14.Sep.B01.06,6,9,0.3,10.4,0.3,36.7,0.5,3435,158.5
14.Sep.B01.07,7,9.1,0.2,10.5,0.2,36.9,0.2,3510,100
```

- There are many ways to do this in Python. We'll start with the simplest (?) and use other methods later

Input from text file

- `s = f.readline()` reads one row from a data file as a single string, here assigned to “`s`”
- can use `t = split(“,”)` to break the string into smaller strings separated by the delineator (here comma, “`,`”), `t` will be of type “list”
- e.g.:

```
In [9]: f = open("/Users/rolandg/Python/cubeData.csv", "r")  
  
In [10]: s = f.readline()  
  
In [11]: print(s)  
14.Sep.B01.01,1,9,0.2,10.4,0.2,36.8,0.3,3400,200  
  
In [12]: t = s.split(",")  
  
In [13]: print(t[0])  
14.Sep.B01.01  
  
In [14]: print(t[9])  
200
```

Homework - I

- If not done yet: Create free account on Slack and join Slack channel
- If not done yet: Install Python framework (e.g., Jupyter) on your laptop
- Send photo with chinese and english name to TA

Homework - II

- Write a short Python script that asks the user for a number and finds the square root of the number using the algorithm on the right
 - do not use any math operators other than $+, -, *, /$
- Send your working Python code to TA
- You can collaborate, but every student needs to write and submit their own code!

- Calculate square root of x , i.e., y such that $y^2 = x$
- Algorithm
 1. guess number y
 2. if y^2 close enough to x : we're done
 3. otherwise average y and x/y to replace y
 4. check if y^2 is close enough to x ; otherwise repeat step 3
- Implementing this algorithm in Python is part 1 of homework...

Homework - III

- Write a short Python script that reads the file “data_HW1.txt”.
The file has 8 columns. Calculate the average of the numbers in each column in a for loop and use `numpy.std()` to calculate the standard deviation of each column (we’ll discuss later how that is defined).
- Send your working Python code to TA
- You can collaborate, but every student needs to write and submit their own code!

x	y	x	y	x	y	x	y
10.0	8.04	10.0	9.14	10.0	7.46	8.0	6.58
8.0	6.95	8.0	8.14	8.0	6.77	8.0	5.76
13.0	7.53	13.0	9.74	13.0	12.74	8.0	7.71
9.0	8.81	9.0	8.77	9.0	7.11	8.0	8.84
11.0	8.33	11.0	9.26	11.0	7.81	8.0	8.47
14.0	9.95	14.0	9.16	14.0	8.84	8.0	7.04
6.0	7.24	6.0	5.13	6.0	6.08	8.0	5.25
4.0	4.25	4.0	3.18	4.0	5.39	19.0	12.50
12.0	10.84	12.0	9.13	12.0	8.15	8.0	5.56
7.0	4.82	7.0	7.26	7.0	6.42	8.0	7.91
5.0	5.63	5.0	4.74	5.0	5.73	8.0	6.89

```
In [16]: x = [1,2,3,4,5,6,7,8,9]
In [17]: np.std(x)
Out[17]: 2.581988897471611
```